



تمرین اول

درخت تصمیم

استاد درس: دکتر آرش عبدی هجران دوست

تدریس یار: سید فرید موسوی

نکات تمرین:

- مهلت تحویل ساعت ۲۳:۵۵ ، ۱۴۰۲/۱۲/۲۹
- مهلت ارسال به هیچ وجه قابل تغییر نیست.
- مواردی که بعد از تاریخ فوق ارسال شوند قابل قبول نبوده و نمره ای نخواهند داشت .
- انجام تمرین تک نفره است. لطفاً به تنهایی انجام شود، در غیر اینصورت نمره منفی در نظر گرفته خواهد شد.
- کل محتوای ارسالی زیپ شود و نام فایل زیپ ارسالی HW1_StudentNumber_YourName باشد.
- محتوای ارسالی دارای راهنما (read me) جهت تسهیل اجرا باشد.
- زبان برنامه نویسی دلخواه است. (پیشنهاد : پایتون)
- در صورت استفاده از زبان پایتون فایل کد ترجیحاً به فرمت ipynb. بوده و فایل کد حتماً به صورت اجرا شده آپلود گردد و از وجود خروجی سلول ها اطمینان حاصل نمایید.
- موارد ارسال شده در تاریخی که بعداً مشخص خواهد شد و متعاقباً اعلام میگردد به صورت آنلاین نیز تحویل گرفته خواهند شد.
- صرفاً آنچه در LMS طبق تاریخ فوق تحویل داده شده است بعداً به صورت حضوری تست شده و توضیح داده میشود.
- تنها تکالیفی که به LMS و قبل از مهلت ارسال، فرستاده می شوند بررسی خواهند شد.
- در صورت داشتن هرگونه سوال میتوانید سوال خود را در گروه تلگرامی درس مطرح کنید .
- حداقل یک ساعت قبل از مهلت ارسال را احتیاطاً هدف قرار دهید، تا مشکلات غیرقابل پیش بینی مانند موارد زیر باعث عدم آپلود پاسخ ها در LMS و ارسال آنها از طریق ایمیل نشوند : (قطعی اینترنت، تنظیم نبودن دقیق ساعت سایت با ساعت گرینویچ، کرش سیستم عامل و نیاز به فرمت، بارش زیبای شهاب سنگ از آسمان و ...)

لطفاً از چت جی پی تی و مدل های هوش مصنوعی مشابه برای پیاده سازی استفاده نکنید، شما خودتان دانشجو

هوش مصنوعی هستید 😊

بخش اول (پیااده‌سازی اولیه)

در ابتدا، می‌خواهیم رده‌بند (classifier) درخت تصمیم را از بیخ و بن (بدون استفاده از کتابخانه‌ی آماده) ، برای داده‌های گسسته، مطابق شبه کد ارائه شده در اسلایدهای کلاس، خودمان پیااده‌سازی نماییم. بهتر است که تاکید کنیم، این بخش بسیار حائز اهمیت میباشد زیرا در ادامه پروژه با کد پیااده سازی شده خودتان قرار است کار کنید.

function DECISION-TREE-LEARNING(*examples, attributes, parent_examples*) **returns**
a tree

if *examples* is empty **then return** PLURALITY-VALUE(*parent_examples*)
else if all *examples* have the same classification **then return** the classification
else if *attributes* is empty **then return** PLURALITY-VALUE(*examples*)
else

$A \leftarrow \operatorname{argmax}_{a \in \text{attributes}} \text{IMPORTANCE}(a, \text{examples})$

tree \leftarrow a new decision tree with root test *A*

for each value v_k of *A* **do**

$\text{exs} \leftarrow \{e : e \in \text{examples} \text{ and } e.A = v_k\}$

subtree \leftarrow DECISION-TREE-LEARNING(*exs, attributes - A, examples*)

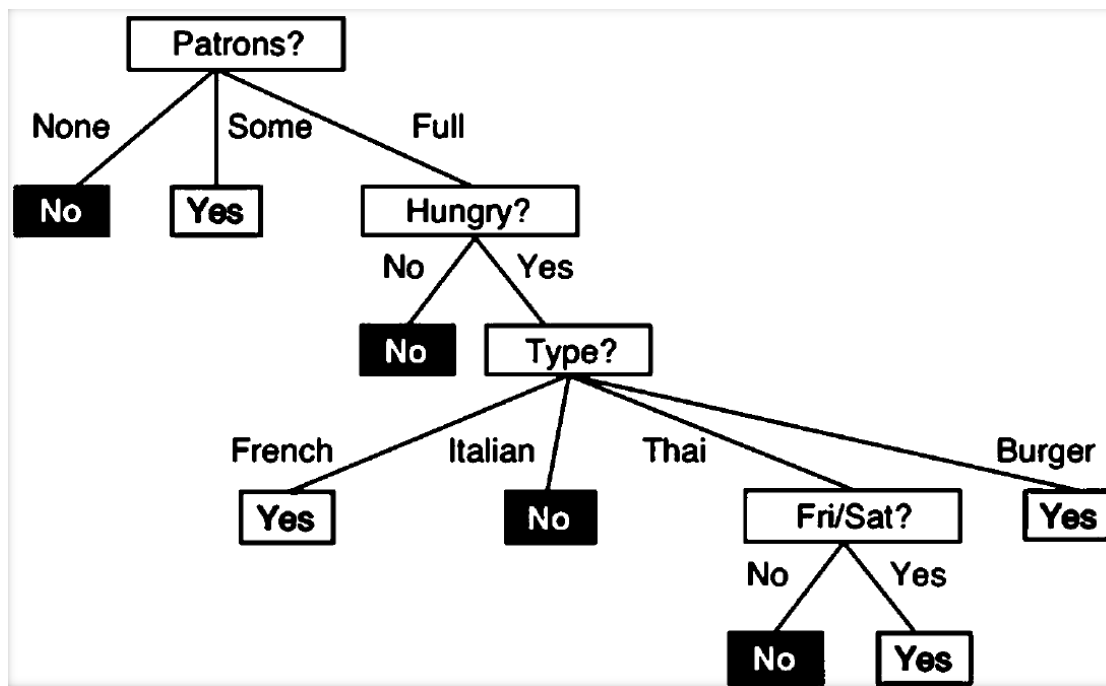
add a branch to *tree* with label ($A = v_k$) and subtree *subtree*

return *tree*

برای اطمینان از صحت پیااده سازی صورت گرفته، میتوانید داده های ۱۲ تایی مثال رستوران (مطرح شده در کلاس درس) را مورد بررسی قرار دهید.

Example	Input Attributes										Goal
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	WillWait
x ₁	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0-10	y ₁ = Yes
x ₂	Yes	No	No	Yes	Full	\$	No	No	Thai	30-60	y ₂ = No
x ₃	No	Yes	No	No	Some	\$	No	No	Burger	0-10	y ₃ = Yes
x ₄	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10-30	y ₄ = Yes
x ₅	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	y ₅ = No
x ₆	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0-10	y ₆ = Yes
x ₇	No	Yes	No	No	None	\$	Yes	No	Burger	0-10	y ₇ = No
x ₈	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0-10	y ₈ = Yes
x ₉	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	y ₉ = No
x ₁₀	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	10-30	y ₁₀ = No
x ₁₁	No	No	No	No	None	\$	No	No	Thai	0-10	y ₁₁ = No
x ₁₂	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30-60	y ₁₂ = Yes

برای این کار تمام ۱۲ داده را به عنوان مجموعه آموزشی در نظر بگیرید (بدون مجموعه‌ی آزمایشی مجزا) و الگوریتم را برای این داده‌ها اجرا نمایید و سعی کنید ساختار درخت آموزش دیده شده را نمایش دهید، خروجی صحیح مطابق تصویر زیر می‌باشد (البته ممکن است دو ویژگی آنتروپی یکسانی داشته باشند و ساختار درخت شما در یک لایه متفاوت باشد).



حتما نیاز نیست درخت به صورت گرافیکی و شماتیک باشد می‌توانید از کتابخانه‌های آماده و یا ساده ترین حالت print درخت استفاده کنید.

بخش دوم (دسته‌بندی)

در این مرحله می‌خواهیم با استفاده از رده‌بند درخت تصمیمی که در مرحله قبل پیاده سازی کرده‌اید، یک مسئله دسته‌بندی با داده‌های واقعی را حل کنیم. برای این کار از مجموعه داده‌گان این [لینک](#) استفاده می‌کنیم که در کنار این فایل با نام Dry_Bean_Dataset.csv نیز آمده است که یک مجموعه داده برای تشخیص دانه های خشک (حبوبات) با استفاده از تعدادی ویژگی می‌باشد. **قابل ذکر هست که مسئله چندکلاسه می‌باشد.**

در ابتدای کار می‌بایست داده‌ها را به دو بخش داده‌های آموزشی (train set) و داده‌های آزمایشی (test set) تقسیم کنید. نحوه بخش بندی داده‌ها به دو بخش آموزش و آزمایش به صورت کاملاً اختیاری و به دلخواه خودتان است (مثلاً ۹۰-۱۰، ۸۰-۲۰ و ...)

همانطور که میدانید ورودی‌های درخت تصمیم باید به صورت گسسته باشد. برای گسسته سازی ورودی‌های از نوع پیوسته (مانند مجموعه داده فعلی) روش‌های مختلفی وجود دارد. ساده‌ترین ایده آن است که برای چنین ویژگی‌هایی، بازه مینیمم تا ماکزیمم اعداد در مجموعه آموزشی را به تعدادی بازه مساوی تقسیم کنید (چه تعداد؟ تعدادهای مختلف را

آزمایش کنید) ایده های بهتر برای گسسته سازی مانند مرتب سازی و انتخاب نقاط برش در هر گره از درخت بر اساس نمونه هایی که در آن گره حاضرند را نیز امتحان کنید. همچنین میتوانید ایده های مطرح شده در کلاس یا ایده های جدید و خلاقانه خود نیز استفاده کنید و نتایج آن را با حالت های قبل (بازه های مساوی یا انتخاب نقاط برش بر حسب مرتب سازی) مقایسه کنید.

سپس با استفاده از الگوریتم نوشته شده درخت تصمیم در بخش قبلی آموزش مدل را بر روی داده های آموزشی انجام دهید (در نظر داشته باشید برای پیاده سازی درخت تصمیم نباید از توابع آماده استفاده کنید، لذا فرمول آنتروپی و ... را باید خودتان پیاده کنید. استفاده از توابع آماده برای قسمتهای بعدی بلامانع است (و حتی توصیه میشود) مثلاً برای خواندن اکسل، نمایش گرافیکی خروجی درخت برای درک شهودی بهتر از فرآیند نحوه تصمیم گیری درخت تصمیم (که الزامی نیست)، نمایش دقت خروجی و ...

همانطور که در کلاس گفته شد معیار هایی برای انتخاب بهترین ویژگی ها وجود دارد، شما با معیار آنتروپی و information gain آشنایی دارید، از آن ها برای feature selection استفاده کنید و نتایج را بررسی کنید، معیار های دیگری نیز در این تسک موجود میباشد برای مثال از معروفترین آن ها میتوان به gini index اشاره کرد. برای بهتر شدن نتیجه خود میتوانید از این معیار استفاده کنید و نتیجه را با دیگر معیار های موجود مقایسه کنید. اگر نتایج متفاوت بود، دلیلی برای این موضوع میتوانید بیان کنید.

برای جلوگیری از بیش برآش (overfitting) از روشهای توقف زودهنگام (early stopping, pre-pruning) برای جلوگیری از برگ های با تعداد نمونه های بسیار کم (یا خیلی gain کوچک) قبل از ساخت کامل درخت و روش های هرس کردن (pruning, post-pruning) پس از آموزش کامل و ساخت درخت، مورد استفاده قرار میگیرد. یک روش برای توقف زودهنگام و یک روش برای هرس کردن را پیاده سازی کنید و با مقادیر مختلف آزمایش کنید. نتایج این دو روش را با یکدیگر مقایسه کنید. دلیلی برای تفاوت این دو پیدا میکنید؟

بعد از فرآیند آموزش درخت تصمیم، مقادیر زیر را برای داده های آموزشی و داده های آزمایشی همهی حالت های گفته شده محاسبه کنید:

۱) صحت (accuracy) ۲) دقت (precision) ۳) فراخوانی (recall) ۴) معیار f1

۵) ماتریس درهم ریختگی (confusion matrix) ۶) نمودار ROC

قابل ذکر هست مقدار شماره ۶ اختیاری میباشد.

دقت کنید که در آزمایشها یک مقدار حداقلی از مقادیر دقت مد نظر است (میتوانید درخت تصمیم نوشته شده با کتابخانه آماده را به عنوان این مقدار مناسب در نظر بگیرید) چنانچه درخت شما مقادیر دقت قابل قبولی نداشت (در مقایسه با درخت آماده موجود در کتابخانه خیلی تفاوت زیاد بود) حتما سعی کنید دقت را افزایش و نتایج را بهبود دهید. (با تغییر مجموعه هایپر پارامترها چه در توزیع داده ها و چه در خود مدل و ...). در غیر این صورت منجر به کسر نمره خواهد شد.

بخش سوم (تقریب تابع)

در این بخش میخواهیم با استفاده از درخت تصمیم تقریب تابع (مسئله رگرسیون) انجام دهیم. برای این کار از مجموعه داده‌ی housing.csv موجود در کنار این فایل، استفاده نماییم که یک مجموعه داده برای تخمین میانگین قیمت یک خانه در کالیفرنیا با استفاده از تعدادی از ویژگی‌های آن خانه میباشد. برای اطلاعات بیشتر درباره این مجموعه داده به این [لینک](#) مراجعه نمایید.

در ابتدای کار نیاز است پس از تقسیم این مجموعه داده به مجموعه آموزشی و آزمایشی، یک پیش پردازش مناسب بر روی آن انجام داده (دقت کنید که مجموعه داده دارای missing value و ... میباشد) و با ویژگی‌های آن بیشتر آشنا شوید و برای ویژگی‌های پیوسته، گسسته سازی به مانند بخش قبلی انجام دهید.

در گام بعدی نیاز است تا درخت نوشته شده در بخش قبلی را با کمی تغییرات (که در کلاس درس به آن اشاره شده است)، به یک درخت تصمیم مسئله رگرسیون تغییر دهید.

در برگ‌ها نیز از روش میانگین‌گیری، میانه‌گیری و fit کردن تابع استفاده نمایید. (نتایج این سه روش را می‌خواهیم در ادامه با یکدیگر مقایسه نماییم)

میتوان مجدداً از هرس کردن برای جلوگیری از بیش برآزش (در صورت لزوم) استفاده کنید.

برای ارزیابی این مدل، معیارهای ارزیابی برای مسئله رگرسیون را مورد مطالعه قرار دهید و هم برای داده‌های آموزشی و هم برای داده‌های آزمایشی این معیارها را محاسبه و تحلیل نمایید.

آنچه تحویل داده میشود:

(۱) کد اجرایی برنامه با توضیحات لازم برای اجرا

(۲) درختی که برای مرحله دوم و سوم پیدا کرده اید (میتوانید گرافیکی نمایش دهید (به هر نحوی که میتوانید) یا به صورت Text با پروتکلی که توضیح میدهید و قابل فهم باشد، بتوان فهمید در هر گره کدام ویژگی با چه مقادیری خروجی تست شده اند و زیر شاخه هایش کدامند.

(۳) گزارش کاملی از مسیر انجام کار، چالش‌هایی که مواجه شده‌اید، اجراهایی که گرفتید و نتایجی که حاصل شده است. گزارش کار از اهمیت بالایی برخوردار است، حجم آن و فرمت استاندارد آن اهمیت ندارد، اما باید نشان دهنده مسیر انجام پروژه، چالش‌ها، راه حل‌ها و نتایج کار شما باشد.